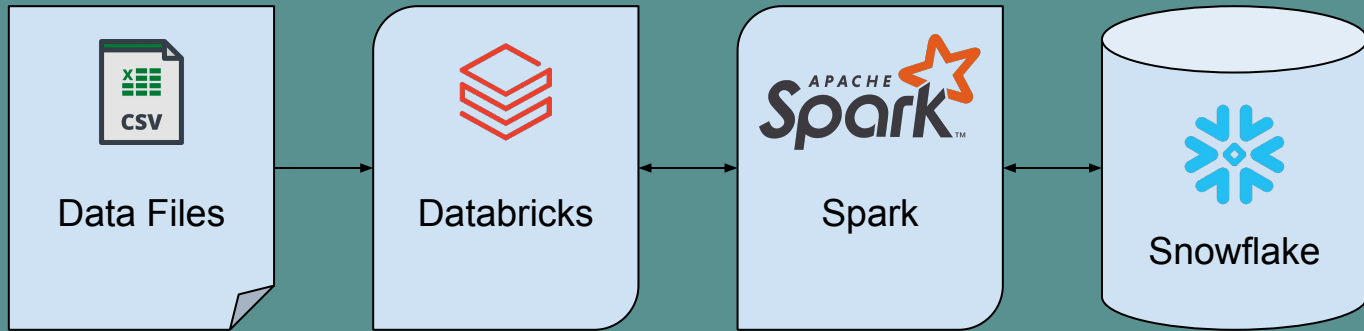# Tools

# Data Flow

# Data

CSV Files

- 14 Airlines, 1 file
- 322 Airports, 1 file
- 3,920,766 Flights, 8 files

3,920,766 Flights!

# Spark
## Extract, Transform, Load (ETL)

Data Science & E... ▾

- Create
- Workspace
- Recents
- Search
- Data
- Compute
- Jobs

ETL (Python)

● My Cluster ▾   File ▾   Edit ▾   View: Standard ▾   Permissions   Run All   Clear ▾   Publish

Read CSV Files
Format Data
Write to Data Warehouse
View Loaded Data

Cmd 1

# Read CSV Files

Cmd 2

```python
def csv_to_df(file_location, file_type):
    # Provide parameters:https://community.cloud.databricks.com/?o=882228783267793#
    #   file location like "/FileStore/tables/phData_challenge/airlines.csv"
    #   file type like "csv"

    # CSV options
    infer_schema = "false"
    first_row_is_header = "true"
    delimiter = ","

    # The applied options are for CSV files. For other file types, these will be ignored.
    df = spark.read.format(file_type) \
        .option("inferSchema", infer_schema) \
        .option("header", first_row_is_header) \
        .option("sep", delimiter) \
        .load(file_location)

    return df

df_airlines = csv_to_df("/FileStore/tables/phData_challenge/airlines.csv", "csv")
df_airports = csv_to_df("/FileStore/tables/phData_challenge/airports.csv", "csv")
df_flights = csv_to_df("/FileStore/tables/phData_challenge/flights/", "csv")
```

# Visualization in Databricks

# Monthly Flights by Airline



Flights by Month

# On Time vs. Late Flights

# Departure & Arrival Delays



Departure & Arrival Delays

# Delay & Cancellation Reasons

# Variety of Routes



Unique Routes

AIRLINE_NAME, AIRLINE_CODE
- Atlantic Southeast Airlines...
- Southwest Airlines Co., WN
- Skywest Airlines Inc., OO
- Delta Air Lines Inc., DL
- American Airlines Inc., AA
- United Air Lines Inc., UA
- American Eagle Airlines I...
- US Airways Inc., US
- JetBlue Airways, B6
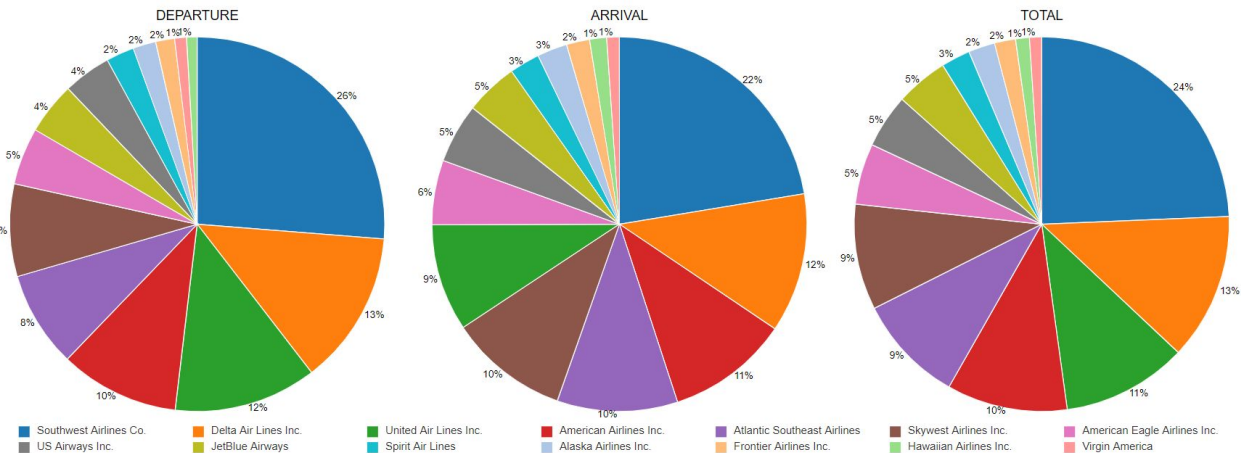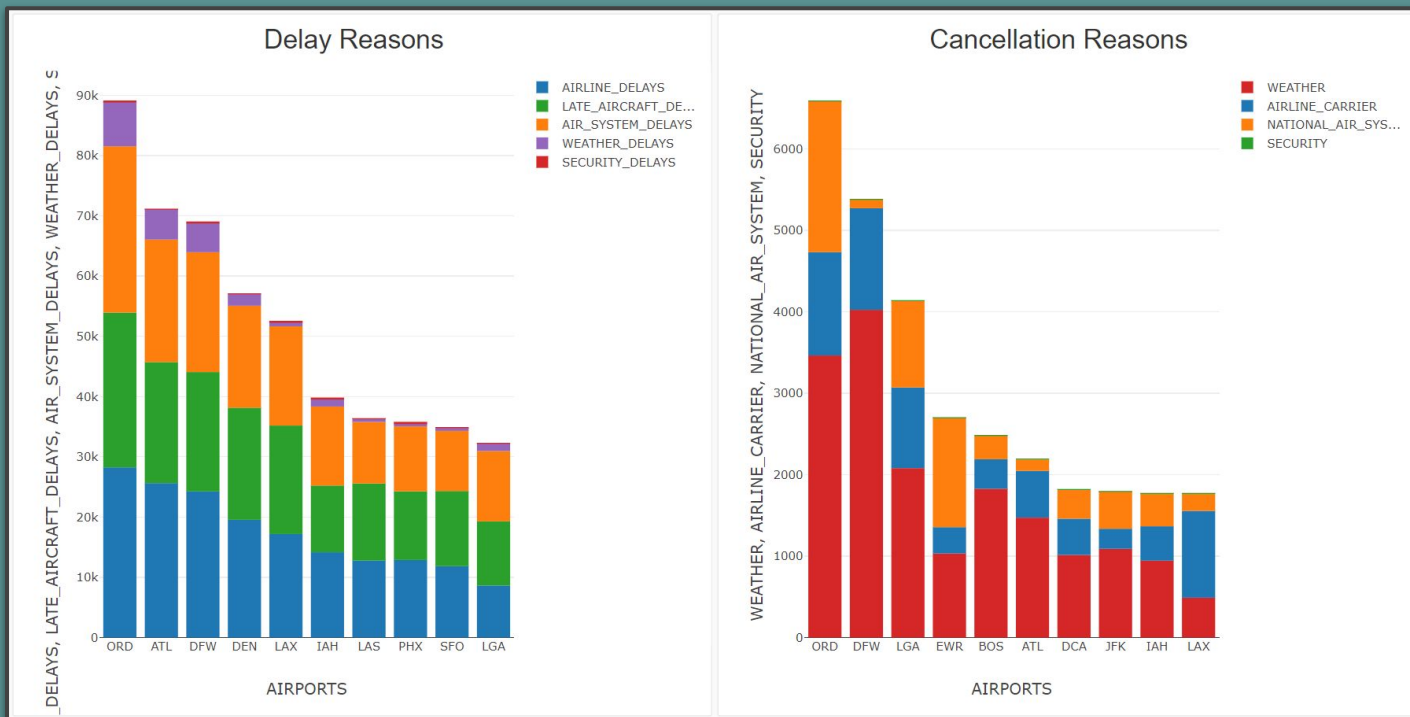- Frontier Airlines Inc., F9
- Spirit Air Lines, NK
- Alaska Airlines Inc., AS
- Virgin America, VX
- Hawaiian Airlines Inc., HA

17%
16%
15%
11%
8%
8%
5%
4%
4%
4%
4%
3%
1% 1%

Prepared by Bogdan Kovch, June 2021
https://github.com/bkovch/DataChallenge
https://www.linkedin.com/in/bogdan-kovch/